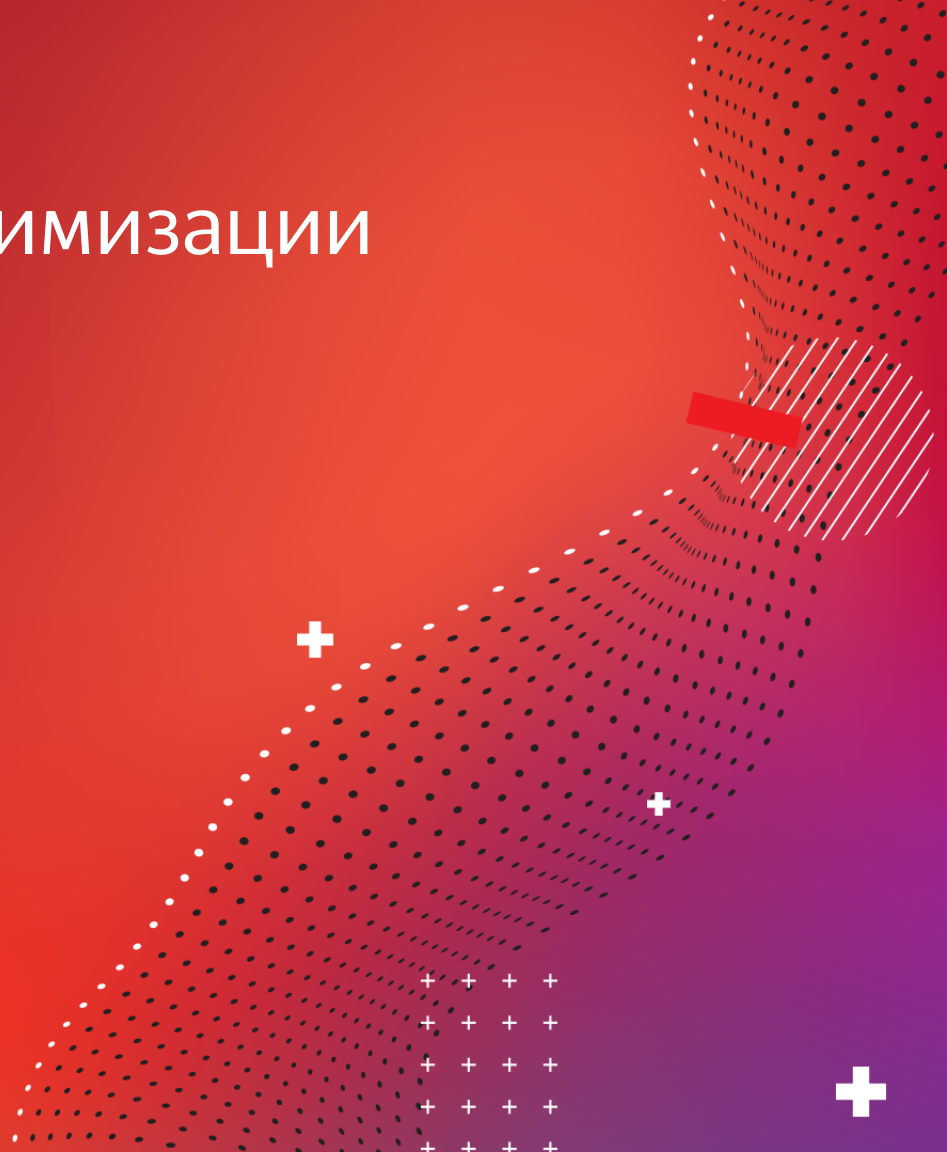


# Необычные случаи оптимизации производительности на примере ClickHouse



**HighLoad++**  
Весна 2021



# Обо мне

Алексей, разработчик ClickHouse.

# ClickHouse

- не тормозил;
- не тормозит;
- будет не тормозить ещё больше.

# Оптимизация производительности

Профилирование на разных нагрузках.

Оптимизация всего, что вылезает.

Про тестирование производительности — смотрите доклад Александра Кузьменкова завтра в 10 утра.



<https://www.techdesignforums.com/practice/technique/winning-at-whac-a-mole-redesigning-an-rf-transceiver/>

# Эпизод 1: MergeTree vs Memory

В ClickHouse есть разные «движки таблиц».

**MergeTree** таблицы хранят данные на диске.

**Memory** таблицы хранят данные в оперативке.

Память быстрее, чем диски\*.

Значит Memory таблицы быстрее, чем MergeTree?

\* Что значит «быстрее»? Скорость последовательного чтения и записи. Задержки случайных чтений и записи. IOPS при заданном параллелизме и распределении нагрузки.

Конечно память может быть медленнее, чем дисковая подсистема, например одноканальная память vs. 10x PCIe 4.0 SSDs.



# MergeTree vs Memory

**Memory** таблицы хранят данные в оперативке.

**MergeTree** таблицы хранят данные на диске, точнее в файловой системе.

Но данные из файловой системы попадают в **page cache**.

И затем читаются уже **из оперативки**.

Значит нет разницы между **Memory** и **MergeTree** таблицами в случае наличия данных в page cache?

# MergeTree vs Memory

Очевидные случаи, когда **MergeTree** быстрее, чем **Memory**.

**MergeTree** таблицы имеют первичный ключ и вторичные индексы, и позволяют читать только нужные диапазоны данных.

**Memory** таблицы позволяют только full scan.

Но этот случай не интересен.

А при full scan может MergeTree быть быстрее, чем Memory?



```

ст 12.168.0.1 | http://sexstant;}p.captcha?p=1&text
| 3567 |
| orton | 2402 | http://tankionliner.by/compani-8-seriildi
| игры лица и гым чан дизайн | http://korer.ru/Session/team/sondakika/gruzki-svjazhene/seasonvali-goda.yandex.do?id=197&img
| url=http | 2166 |
| imgsrc | 1848 | http://haberi.hurriyet.com/bets_z_krymskih-porno/438828f3c987176434229589%2FarchPeriod=60&wp
=1&text= | 1659 | http://auto.ru/chatay-barana.ru/traction.html#maybettaya
| брызговики | 1549 | http://korer.ru/categories.ru/?vkb
| индийский афтозный | 1480 | http://korer.ru/categories.ru/?vkb
| ооооотводка и | 1247 | http://hurpass
| выступная мужчин | 1112 | http://images/inbox/?_lld
| юность |

```

10 rows in set. Elapsed: 0.221 sec. Processed 8.87 million rows, 880.65 MB (40.15 million rows/s., 3.98 GB/s.)

```

milovidov-desktop :) Bye.
milovidov@milovidov-desktop:~/work/clickhouse-presentations/highload2021/pictures$
milovidov@milovidov-desktop:~$ clickhouse-client
ClickHouse client version 21.6.1.1.
Connecting to localhost:9000 as user default.
Connected to ClickHouse server version 21.6.1 revision 54448.

```

```
milovidov-desktop :) 
```

```

~: mc highload2021: mc highload2018: mc ~: bash ~: bash ~: bash ~: bash ~: sudo openmpn

```

index.html

```

1 | <!-- MergeTree</b> таблицы имеют первичный ключ и вторичные индексы,<br/>и позволяют читать только нужные диапазоны данных.</p>
2 | <!-- MergeTree</b> таблицы позволяют только full scan.</p>
3 |
4 | <p style="color: gray;">Жо этот случай не интересен.</p>
5 |
6 | <p style="color: green;">А при full scan может MergeTree быть быстрее, чем Memory?</p>
7 | </section>
8 |
9 | <!-- clickhouse-benchmark <<< SELECT SearchPhrase, any(replaceAll(substring(URL, 1, 100), 'google', 'yandex')) AS s, count() AS c FROM test.hits GROUP BY SearchPhrase ORDER BY c DESC LIMIT 10 -->
10 |
11 | <section class="slide">
12 |   <video style="width: 100%;"><source src="video/memory_is_slower.ogv" type="video/ogg"></video>
13 | </section>
14 |

```

Find: #(<id>)

Replace: <a href="https://github.com/ClickHouse/ClickHouse/pull/11">#11</a>

Mode: Regular expression

Replace Replace All Find All

ов HL++ 2021 весна Moscow

уже на финишной прямой перед конференцией! Главн...

Sign up

milicheva

иве!

admin

на финишной прямой перед конференцией!

поинты:

ти для докладчиков, 16 мая с 18:00 до 21:00 в ресторане "мус" по адресу: Большая Дорогомиловская, 4

тельно пришлите, сегодня вашу презентацию на [shighload@ontico.ru](mailto:shighload@ontico.ru). Финальную версию презентации необходимо прислать)

вы хотите показывать свою презентацию со своего лаптопа, нужно нам сообщить.

ас будут какие-либо вопросы, пишите или звоните, мы на площадке. Юм - 8 905 012 40 80 - 8 913 743 20 44

адур! ранее пропустили в общем чате более полную информацию по основным организационным вопросам, все дублирую, чтобы не потерялось) Прочитайте, это будет полезно

таци:

пуста, убедиться, что в вашей презентации нет ошибок с авторскими правами. А то YouTube ругается до блокирования видео, если есть - просим заменить, или убрать звук. Кстати, "чужие" картинки тоже можно не использовать, YouTube еще не научился их авать, но лучше не рисковать

Будьте, пожалуйста, прислать презентацию на вычитку [shighload@ontico.ru](mailto:shighload@ontico.ru). Кто уже прислал, мы получили, ждем, редактор скоро ответит.

еское оснащение во время выступления:

ик, с которого запускается презентация, будет стоять на 3-м этаже в вестибюле доклада (хорошо оборудован). На экране сзади спикера выводится презентация, не дублирует (расположен перед спикером) - также текущий слайд, и вы можете показывать презентацию со своего ноутбука, только сообщите мне об этом.

ление: до и после

зидки для вас будет выделена специальная стойка с указателем "Докладчик", где вы сможете найти бейдж спикера, минуя очередь. Комната докладчиков - в которой вы сможете оставить пакеты, отдохнуть, подготовиться к выступлению, выступлению, просим вас подойти в зал за 15 минут до начала, нужно успеть проверить финально презентацию, настроить микрофон и провести для вас инструктаж. Презентацию можно проверить заранее во время репетиции в зале.

выступления мы предлагаем использовать головные телефоны.

ть докладов — по 30-40 минут + 10-15 минут ответы на вопросы.

да в течение 30 минут мы просим вас находиться в зоне у выхода из зала, где все желающие смогут задать вопросы. Также мы предлагаем использовать цифровые кулеры, участники из

# MergeTree vs Memory

Неочевидные случаи, когда **MergeTree** быстрее, чем **Memory**.

**MergeTree** таблицы хранят данные в отсортированном порядке по первичному ключу.

Некоторые алгоритмы в ClickHouse эксплуатируют преимущества локальности данных, если она есть (fast path).

Например, если при GROUP BY подряд дважды встретилось одно и то же значение, то мы не делаем повторный поиск в хэш-таблице.

Про хэш-таблицы в ClickHouse смотрите доклад Максима Киты завтра в 12:50.

А если данные в таблицах находятся в одинаковом порядке, может ли MergeTree быть быстрее, чем Memory?

# Как обрабатываются данные в ClickHouse

Данные в ClickHouse хранятся по столбцам  
и обрабатываются тоже по столбцам.

## Array of Structures

```
struct Point3d
{
    float x;
    float y;
    float z;
};
std::vector<Point3d> points;
```

## Structure of Arrays

```
struct Points
{
    std::vector<float> x;
    std::vector<float> y;
    std::vector<float> z;
};
```

# Как обрабатываются данные в ClickHouse

Данные в ClickHouse хранятся по столбцам  
и обрабатываются тоже по столбцам. **По кусочкам столбцов.**

```
struct Chunk
{
    std::vector<float> x;
    std::vector<float> y;
    std::vector<float> z;
};
```

**`std::vector<Chunk> chunks;`**

— Morsel-based processing.

# Как именно читаются данные из таблицы?

В случае **MergeTree**:

- читаем сжатые файлы из файловой системы;
- вычисляем и сверяем чексуммы;
- разжимаем сжатые блоки;
- десериализуем кусочки столбцов;
- обрабатываем их;



**Это не оптимально**

# Как именно читаются данные из таблицы?



В случае **Memory**:

- в оперативке уже находятся готовые кусочки столбцов,  
обрабатываем их;

**Это более оптимально**

# Что именно происходит при чтении?

В случае **MergeTree**:

1. Читаем сжатые файлы из файловой системы:

- читать можно с помощью синхронного (**read/pread, mmap**) или асинхронного (**AIO, uring**) ввода-вывода;
- в случае синхронного ввода-вывода, можно использовать (обычный **read** или **mmap**) или не использовать **page cache (O\_DIRECT)**;
- если читать из **page cache** без **mmap**, то будет **копирование** из **page cache** в **userspace**;
- читаем сжатые данные — если коэффициент сжатия большой, то доля времени в обработке запроса маленькая;

# Что именно происходит при чтении?

В случае **MergeTree**:

2. Разжимаем сжатые блоки:

- по-умолчанию используется LZ4\*;
- можно выбрать как более сильный метод сжатия (ZSTD), так и более слабый, например вообще без сжатия (NONE);
- иногда NONE внезапно работает медленнее, с чего бы это?
- а блоками какого размера были сжаты данные?  
и как это влияет на скорость?

\* Смотрите доклад «Как ускорить разжатие LZ4» с HighLoad++ Siberia 2018.



# Что именно происходит при чтении?

В случае **MergeTree**:

3. Десериализуем кусочки столбцов:

- десериализации как таковой нет;
- это просто перекладывание данных (memcpy);
- а зачем вообще нужно перекладывать данные?

# Отличие MergeTree и Memory

В случае **Memory**:

— готовые кусочки столбцов в оперативке.

В случае **MergeTree**:

— кусочки столбцов формируются динамически при чтении.

**MergeTree** делает больше работы,  
но может ли это иногда быть оптимальнее?

# MergeTree vs Memory

В случае **MergeTree**:

- кусочки столбцов формируются динамически при чтении, и их размер в числе строк может выбираться адаптивно для **кэш-локальности**!

# Кэш-локальность

**С какой скоростью работает оперативка?**

— какая оперативка, на какой машине?

**С какой скоростью работает кэш?**

— кэш какого уровня, на каком CPU?

— один или все вместе?

**С какой скоростью чего?**

— throughput, latency?..

# Эпизод 2: сжатие данных тормозит?

В ClickHouse данные по-умолчанию хранятся сжатыми.

При записи сжимаются, при чтении — разжимаются.

Профилируем запросы...

В топе по CPU — функция `LZ4_decompress_safe`.

🤔 Чтобы всё ускорить, надо просто убрать сжатие данных?



# Пробуем убрать сжатие данных

Но ничего хорошего из этого не выходит:

1. Убрали сжатие данных и теперь они не помещаются на диск.
2. Убрали сжатие данных и теперь чтение с диска тормозит.
3. Убрали сжатие данных и теперь меньше данных помещается в page cache.

...

Но даже если несжатые данные помещаются целиком в оперативку — **имеет ли смысл не сжимать** их?



# Что быстрее: разжатие или memcpu?

Функцию **memcpu** используют как baseline самого слабого сжатия или разжатия в бенчмарках.

Конечно, это самый быстрый эталон для сравнения.

Пример:

- memcpu: **12 ГБ** в секунду.
- LZ4 decompression: **2..4 ГБ** разжатых данных в секунду.

Вывод: memcpu быстрее, чем разжатие LZ4?



# Что быстрее: разжатие или memcached?

Рассмотрим сценарий:

- данные хранятся в оперативке;
- данные обрабатываются по блокам;
- каждый блок достаточно небольшой и помещается в кэш CPU;
- обработка каждого блока помещается в кэш CPU;
- **данные обрабатываются в несколько потоков;**

Данные читаются из оперативки, дальше используется только кэш CPU.

# Что быстрее: разжатие или memsru?

Пример: Ryzen 3950 (16 ядер)

- memsru:  $16 \times 12 \text{ ГБ} = 192 \text{ ГБ}$  в секунду.
- LZ4 decompression:  $16 \times 2..4 \text{ ГБ} = 32..48 \text{ ГБ}$  разжатых данных в секунду.
- скорость чтения из памяти: **30 ГБ\*** в секунду.

В случае memsru чтение упирается в скорость памяти.

Но если используется сжатие, то из памяти читается меньше данных.  
**Память работает как диск.** Разжатие LZ4 быстрее, чем memsru?

\* память двухканальная, но работает не на максимальной частоте.  
По спецификации для этого CPU до 48 ГБ в секунду.

# Что быстрее: разжатие или тетсру?

Пример: 2 × AMD EPYC 7742 (**128 ядер**)

8 channel memory, max throughput **190 GiB/s**

Для этого сервера работа с данными,  
сжатыми LZ4, также будет быстрее.

Но если ядер меньше — уже не всё однозначно.

Если данные хорошо сжаты, то разжатие всё-таки упирается в CPU,  
а значит, его можно ускорить!

# Оптимизации в ClickHouse

Для **Memory** таблиц:

- Уменьшили размер блока при записи для лучшей кэш-локальности обработки данных [#20169](#).
- Возможность сжатия Memory таблиц [#20168](#).

Для **MergeTree** таблиц:

- Убрали лишнее копирование для режима сжатия NONE [#22145](#).
- Возможность отключить чексуммы при чтении [#19588](#), но использовать эту возможность не надо.
- Возможность чтения с помощью mmap [#8520](#), чтобы убрать лишнее копирование из page cache а также кэш memory mappings [#22206](#).

```
~ : mc — Konsole
File Edit View Bookmarks Settings Help

80.000%      0.218 sec.
90.000%      0.219 sec.
95.000%      0.219 sec.
99.000%      0.219 sec.
99.900%      0.219 sec.
99.990%      0.219 sec.

^CStopping launch of queries. SIGINT received.

Queries executed: 37.

localhost:9000, queries 37, QPS: 4.631, RPS: 41091849.686, MiB/s: 3889.060, result RPS: 46.306, result MiB/s: 0.004.

0.000%      0.208 sec.
10.000%     0.211 sec.
20.000%     0.211 sec.
30.000%     0.212 sec.
40.000%     0.214 sec.
50.000%     0.216 sec.
60.000%     0.216 sec.
70.000%     0.217 sec.
80.000%     0.219 sec.
90.000%     0.223 sec.
95.000%     0.225 sec.
99.000%     0.234 sec.
99.900%     0.234 sec.
99.990%     0.234 sec.

milovidov@milovidov-desktop: ~$
```



# Выводы

Чтобы оптимизировать производительность, нужно всего лишь:

- точно знать, что делает ваш код;
- профилировать систему на реалистичных сценариях нагрузки;
- представлять возможности железа;
- ...
- не забывать что в системе много ядер, а у процессора есть кэш;  
не путать latency и throughput :)

# Спасибо!